

# Approximation Algorithms for the Shortest Common Superstring Problem

JONATHAN S. TURNER  
 Washington University  
 St. Louis, Missouri  
 jst@wucs.UUCP

## Abstract

The object of the shortest common superstring problem (SCS) is to find the shortest possible string that contains every string in a given set as substrings. As the problem is NP-complete, approximation algorithms are of interest. The value of an approximate solution to SCS is normally taken to be its length, and we seek algorithms that make the length as small as possible. A different measure is given by the sum of the overlaps between consecutive strings in a candidate solution. When considering this measure, the object is to find solutions that make it as large as possible. These two measures offer different ways of viewing the problem. While the two viewpoints are equivalent with respect to optimal solutions, they differ with respect to approximate solutions. We describe several approximation algorithms that produce solutions that are always within a factor of two of optimum with respect to the overlap measure. We also describe an efficient implementation of one of these, using McCreight's compact suffix tree construction algorithm. The worst-case running time is  $O(m+n \log n)$  for small alphabets, where  $m$  is the sum of the lengths of all the strings in the set and  $n$  is the number of strings. For large alphabets, the algorithm can be implemented in  $O(m \log m)$  time by using Sleator and Tarjan's lexicographic splay tree data structure.

## 1. Introduction

Let  $s_1 = a_1 \dots a_r$  and  $s_2 = b_1 \dots b_s$  be strings over some finite alphabet  $\Sigma$ . We say that  $s_1$  is a *substring* of  $s_2$  if there is an integer  $i \in [0, s-r]$  such that  $a_j = b_{i+j}$  for  $1 \leq j \leq r$ . We also say in this case that  $s_2$  is a *superstring* of  $s_1$ .

An instance of the *shortest common superstring problem* (SCS) is a set of strings  $S = \{s_1, \dots, s_n\}$  over a finite alphabet  $\Sigma$ . The object of the problem is to find a minimum length string that is a superstring of every  $s_i \in S$ . We let  $\phi^*(S)$  denote the length of a minimum length superstring.

We say that a set of strings is *substring free* if no string in the set is a substring of any other. We will generally limit our attention to substring free sets. This involves no loss of generality, since any set of strings has a unique substring free subset which has the same solutions as the original set.

We have presented the problem in the conventional way, with the object being to minimize the solution length. It is useful to consider an alternative viewpoint as well. One can view the object of the problem as being to find an ordering of the strings that maximizes the amount of overlap between consecutive strings. To make this precise we need a few definitions.

Let  $s_1 = a_1 \dots a_r$  and  $s_2 = b_1 \dots b_s$  be strings. We define

$$\psi(s_1, s_2) = \max \{k \geq 0 \mid a_{r-k+i} = b_i \ 1 \leq i \leq k\}$$

If  $\psi(s_1, s_2) = k$  then  $s_1 \circ s_2$  is defined to be the string  $a_1 \dots a_r b_{k+1} \dots b_s$ . We note that if  $s_1, s_2, s_3$  are strings, none of which is a substring of another, then  $s_1 \circ (s_2 \circ s_3) = (s_1 \circ s_2) \circ s_3$ ; that is, the

---

The research described here was supported by the National Science Foundation (grant #DCR-8409435) and the National Institutes of Health (grant #RR-01380)

overlapping operation is associative for substring free sets. Consequently, we may write  $s_1 \circ s_2 \circ \dots \circ s_n$  with no ambiguity.

Let  $\pi$  be a permutation on  $\{1, \dots, n\}$ . We will usually write  $\pi_i$  for  $\pi(i)$ . We define

$$\psi_\pi(s_1, \dots, s_n) = \sum_{i=1}^{n-1} \psi(s_{\pi_i}, s_{\pi_{i+1}})$$

and  $\phi_\pi(s_1, \dots, s_n) = |s_{\pi_1} \circ \dots \circ s_{\pi_n}|$ . Note that for any instance  $S = \{s_1, \dots, s_n\}$  of SCS,

$$\phi_\pi(S) = \|S\| - \psi_\pi(S)$$

where  $\|S\| = \sum_{i=1}^n |s_i|$ . In particular,

$$\phi^*(S) = \|S\| - \psi^*(S) \quad \text{where} \quad \psi^*(s_1, \dots, s_n) = \max_{\pi} \psi_\pi(s_1, \dots, s_n)$$

Hence, we can view the object of the SCS problem as being to find a mapping  $\pi$  that maximizes  $\psi_\pi$ . Let  $A$  be an algorithm for SCS which given a collection of strings  $S = \{s_1, \dots, s_n\}$  produces a mapping  $\pi = \pi_A(S)$ . We define  $\psi_A(S) = \psi_{\pi_A}(S)$  and  $\phi_A(S) = \phi_{\pi_A}(S)$ .

SCS was shown to be NP-complete by Maier and Storer in [9]. Another, and more elegant proof appears in [4] and [5]. One obvious application for the problem is data compression. Storer and Szymanski [16] for example, consider a fairly general model of data compression which includes SCS as an important special case. See also [10]. Another application is to DNA sequencing. SCS is one of the simplest models for the problem of recovering DNA sequencing information from experimental data [6,13,15]. To our knowledge the only approximation algorithm to be discussed in the literature is a simple greedy algorithm which is treated briefly by Gallant in [5]. Gallant claims that for this algorithm, which we refer to as SGREEDY,  $\phi_{\text{SGREEDY}}(S) \leq (3/2)\phi^*(S)$  for all collections of strings  $S$ . We show that this is not in fact true by displaying a set of strings  $S$  for which  $\phi_{\text{SGREEDY}}(S) \approx 2\phi^*(S)$ . We have found no worse example problem than this, but have also been unsuccessful in proving an upper bound on the performance of this algorithm in terms of the length measure. On the other hand, we do show that  $\psi^*(S) \leq 2\psi_{\text{SGREEDY}}(S)$ .

In section 2 we relate SCS to the longest path problem (LPP) in graphs by describing a transformation from SCS to LPP that preserves solution values with respect to the overlap measure. We then construct three approximation algorithms for LPP, two based on matching and the third a greedy heuristic. By virtue of the transformation from SCS, all three are also approximation algorithms for SCS. We show that the greedy heuristic for LPP always produces solutions within a factor of three of the optimum value. In section 3, we show that the instances of LPP that result from our transformation from SCS have a special structure that allows us to obtain a tighter bound. We also describe an efficient implementation of this greedy algorithm for strings using a compact representation of suffix trees. In section 4, we relate SCS to the traveling salesman problem (TSP) by another transformation that preserves solution values, this time with respect to the length measure. The instances of TSP arising from this transformation are asymmetric, but satisfy the triangle inequality. There are no approximation algorithms known for this problem with provably good worst-case performance, nor have we succeeded in finding any. Nevertheless, this transformation means that if such an algorithm is found, it can be used for SCS as well as TSP. If on the other hand, it turns out that approximating this version of TSP is hard, then any approximation algorithm for SCS, will have to make use of special structural properties present in the instances of TSP that arise from this transformation.

## 2. SCS and the Longest Path Set Problem

In this section we relate SCS to the *longest path problem* (LPP) in graphs. An instance of the longest path problem is a complete directed graph  $G = (V, E)$  with each edge  $(u, v)$  having a non-negative integer length  $\ell(u, v)$ . The *length of a path*  $p$  in  $G$  is defined to be the sum of the lengths of its edges and is denoted  $\lambda_p(G, \ell)$ . The object of the longest path problem is to find a Hamiltonian path  $p$  (that is a path including every vertex) in  $G$  that maximizes  $\lambda_p(G, \ell)$ . The length of such a longest path is denoted  $\lambda^*(G, \ell)$ .

Let  $S = \{s_1, \dots, s_n\}$  be an instance of SCS. We define LPP( $S$ ) to be an instance  $(G, \ell)$  of LPP with

$$\begin{aligned} V &= \{u_1, \dots, u_n\} & E &= V \times V \\ \ell(u_i, u_j) &= \psi(s_i, s_j) & 1 \leq i, j \leq n, i \neq j \end{aligned}$$

```

function edgeset MATCH(digraph  $G = (V, E)$ , edgelengths  $\ell$ )
  edgeset  $P, M$ ;
   $P := \emptyset$ ;
  do  $E \neq \emptyset \rightarrow$ 
     $M := \text{MAXMATCH}(G, \ell)$ ;
     $P := P \cup M$ ;
    for  $(u, v) \in M \rightarrow$ 
      Delete from  $G$ , all edges of the form  $(u, x)$  or  $(y, v)$ ;
      Collapse  $u$  and  $v$  into a single vertex;
    rof;
  od;
  return  $P$ ;
end;

```

Figure 1: Matching Algorithm for LPP

Let  $\pi$  be a permutation on  $\{1, \dots, n\}$ . We can view  $\pi$  as defining a Hamiltonian path  $u_{\pi_1}, \dots, u_{\pi_n}$  in  $G$ . We let  $\lambda_\pi(G, \ell)$  denote the length of this path. We now state a trivial, but useful theorem.

**THEOREM 1.** *Let  $S = \{s_1, \dots, s_n\}$  be an instance of SCS,  $(G, \ell) = \text{LPP}(S)$  and let  $\pi$  be any permutation on  $\{1, \dots, n\}$ .  $\lambda_\pi(G, \ell) = \psi_\pi(S)$ . In particular,  $\lambda^*(G, \ell) = \psi^*(S)$ .*

The theorem implies that any approximation algorithm for LPP is an approximation algorithm for SCS with respect to the overlap measure. In the remainder of this section, we present three simple approximation algorithms for LPP.

**Matching Algorithm.** A *matching* in a graph  $G = (V, E)$  is a set of edges, no two of which share a common vertex. A maximum matching in a graph with edge lengths  $\ell(e)$  is a matching  $M$  that maximizes  $\ell(M)$ . We define  $\mu^*(G, \ell) = \max_M \ell(M)$  to be the value of a maximum matching. There are algorithms for finding maximum matchings having running times of  $O(n^3)$  (where  $n = |V|$ ) [17].

Our first algorithm for LPP is based on the observation that any matching for an instance  $(G, \ell)$  of LPP can be extended to a path (since  $G$  is assumed to be complete) and a maximum matching must have total length at least half that of a longest path. (Recall that we are restricting attention to non-negative weights.)

**THEOREM 2.** *If  $(G = (V, E), \ell)$  is an instance of LPP then  $\lambda^*(G, \ell) \leq 2\mu^*(G, \ell)$ .*

Theorem 2 provides the basis for our first approximation algorithm shown in Figure 1. The procedure MATCH starts by finding a maximum matching in  $G$ , then removes edges that are ruled out by the selected edges, collapses the selected edges into single vertices and then repeats the process on the new graph. To see that the algorithm does construct a Hamiltonian path, note the following: (1) the edge eliminations ensure that the set  $P$  never contains two edges leaving a common vertex or entering a common vertex, (2) the collapsing of edges into single vertices prevents creation of cycles and (3) since the original graph is assumed to be complete, the algorithm will halt only when a complete Hamiltonian path has been constructed.

Theorem 2 implies that  $\psi^*(G, \ell) \leq 2\psi_{\text{MATCH}}(G, \ell)$  for any instance  $(G, \ell)$  of LPP. The running time of MATCH is determined primarily by the matching algorithm used. Assuming a matching algorithm that runs in  $O(n^3)$  time, we get a running time of  $O(n^3 \log n)$  for MATCH.

**Directed Matching Algorithm.** A *directed matching* in a digraph  $G = (V, E)$  is a set of edges, no two of which enter a common vertex and no two of which leave a common vertex. In other words, it is a subgraph of  $G$  comprising a collection of disjoint paths and cycles. A maximum directed matching in a graph  $G$  with edge lengths  $\ell(e)$  is a directed matching  $M$  that maximizes  $\ell(M)$ . We define  $\delta^*(G, \ell) = \max_M \ell(M)$  to be the value of a maximum directed matching (where in this case,  $M$  ranges over all directed matchings of  $G$ ). There are algorithms for finding maximum directed matchings with running times of  $O(n^{5/2})$  [17].

Given any matching  $M$ , let  $M^-$  be a subset of  $M$  obtained by discarding a least cost edge from each cycle in  $M$ . Our next algorithm for LPP is based on the observation contained in the next theorem.

```

function edgeset DIMATCH(digraph  $G = (V, E)$ , edgelengths  $\ell$ )
  edgeset  $P, M$ ;
   $P := \emptyset$ ;
  do  $E \neq \emptyset \rightarrow$ 
     $M := \text{MAXDIMATCH}(G, \ell)$ ;
     $M^- := M - \text{one least cost edge from each cycle of } M$ ;
     $P := P \cup M^-$ ;
    for each path  $(u_1, \dots, u_r) \in M^- \rightarrow$ 
      Delete from  $G$ , all edges of the form  $(u_i, x)$ ,  $2 \leq i \leq r$ ;
      Delete from  $G$ , all edges of the form  $(x, u_i)$ ,  $1 \leq i \leq r - 1$ ;
      Delete from  $G$  the edge  $(u_r, u_1)$ , if present;
      Collapse the path into a single vertex;
    rof;
  od;
  return  $P$ ;
end;

```

Figure 2: Directed Matching Algorithm for LPP

**THEOREM 3.** *Let  $(G = (V, E), \ell)$  be an instance of LPP, let  $M$  be a maximum directed matching of  $G$  and let  $k$  be the minimum number of edges in any cycle defined by  $M$ .  $\lambda^*(G, \ell) \leq \frac{k}{k-1} \ell(M^-)$ . In particular,  $\lambda^*(G, \ell) \leq 2\ell(M^-)$ .*

We note that  $\delta^*(G, \ell) \geq \lambda^*(G, \ell)$ . Hence, it provides a measure of how close a given solution is to optimal. We expect that the solutions obtained by breaking cycles will often be much closer to optimal than the bound in the theorem implies.

Theorem 3 provides the basis for our next approximation algorithm for LPP, shown in Figure 2. This algorithm constructs a maximum directed matching  $M$  in  $G$ , then breaks all the cycles in  $M$  and constructs a new graph in which the paths of  $M$  correspond to vertices. It then proceeds by finding a maximum directed matching in the new graph, continuing in this fashion until a Hamiltonian path in the original graph has been found.

Theorem 3 implies that  $\psi^*(G, \ell) \leq 2\psi_{\text{DIMATCH}}(G, \ell)$  for any instance  $(G, \ell)$  of LPP. The running time of DIMATCH is determined primarily by the directed matching algorithm used. Assuming an algorithm that runs in  $O(n^{5/2})$  time, we get a running time of  $O(n^{5/2} \log n)$  for DIMATCH.

DIMATCH is essentially an adaptation of an algorithm for the asymmetric traveling salesman problem (TSP) described by Karp in [7]. Karp's algorithm has poor worst-case performance for TSP, but performs well in a probabilistic sense for instances in which inter-city distances are selected uniformly on the interval  $[0, 1]$ . We have simply adapted his algorithm to the longest path problem (simplifying it slightly in the process), and observed that its worst-case performance is provably good in this context.

**Greedy algorithm.** The algorithms considered above are both fairly complicated and time consuming because they require the calculation of maximum weighted matchings. Another algorithm that is worth considering is the simple greedy algorithm that scans the edges in non-increasing order of length and selects an edge  $(u, v)$  if it has not previously selected an edge of the form  $(u, x)$  or  $(y, v)$  and if the collection of paths constructed so far does not include a path from  $v$  to  $u$ . The next theorem gives a worst-case bound on the performance of the greedy algorithm.

**THEOREM 4.** *If  $(G, \ell)$  is an instance of LPP then  $\lambda^*(G, \ell) \leq 3\lambda_{\text{GREEDY}}(G, \ell)$ .*

*proof.* Let  $F$  be the set of edges in some optimum solution to  $(G, \ell)$ . Let  $H = \{h_1, \dots, h_s\}$  be the set of edges chosen by the greedy algorithm in the order in which they were selected (that is,  $h_1$  was selected first,  $h_2$  second, and so forth).

We say an edge is *permissible* at some stage of the execution of the algorithm if its selection hasn't been precluded by earlier selections. Define  $H_i$  to be the set of edges which are permissible just before  $h_i$  is selected, but not permissible after  $h_i$  is selected.

Let  $h_i = (u, v)$  and consider the situation just before  $h_i$  is selected by the greedy algorithm. At this point,  $u$  is the last vertex of some path constructed by the algorithm and  $v$  is the first vertex of

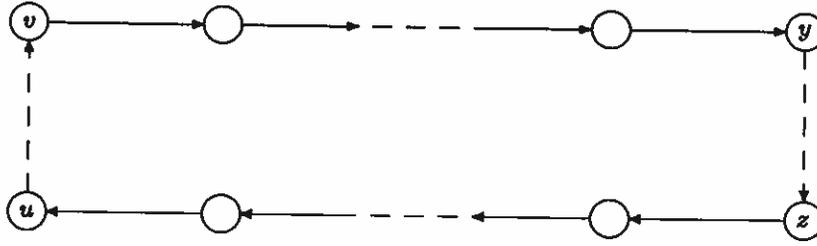


Figure 3: Illustration for Theorem 2.4

```

procedure PGREEDY(digraph  $G = (V, E)$ , edgelengths  $\ell$ ,
mapping  $left, right : V \mapsto V \cup \{\text{null}\}$ )
vertex  $u, v$ ; mapping  $leftend, rightend : V \mapsto V$ ;
for  $u \in V \rightarrow$ 
     $left(u), right(u) := \text{null}$ ;
     $leftend(u), rightend(u) := u$ ;
rof;
Sort  $E$  from longest to shortest;
for  $(u, v) \in E \rightarrow$ 
    if  $right(u) = \text{null}$  and  $left(v) = \text{null}$  and  $v \neq leftend(u)$ 
         $right(u), left(v) := v, u$ ;
         $rightend(leftend(u)) := rightend(v)$ ;
         $leftend(rightend(v)) := leftend(u)$ ;
    fi;
rof;
end

```

Figure 4: Greedy Algorithm for LPP

some path (one or both paths may contain just a single vertex). Let  $z$  be the first vertex on the path containing  $u$  and let  $y$  be the last vertex on the path containing  $v$  as shown in Figure 3.

If there is an edge from  $y$  to  $z$  that is permissible before the selection of  $h_i$  then that edge is a member of  $H_i$ . All other members of  $H_i$  have the form  $(u, x)$  or the form  $(x, v)$ .

Next, note that  $\ell(h_i) = \max \{\ell(e) \mid e \in H_i\}$  and that  $(H_1, \dots, H_s)$  is a partition of  $E$ . Finally note that for  $i \in [1, s]$ ,  $|F \cap H_i| \leq 3$ . This yields the theorem.  $\square$

Figure 4 is a sketch of an implementation of the greedy algorithm. Upon return, the mappings  $left$  and  $right$  give the left and right neighbors of each vertex in the solution path. If  $left(u)$  is null, then  $rightend(u)$  gives the vertex at the end of the path containing vertex  $u$  in the current partial solution.  $leftend(u)$  is similar. The running time for this implementation is  $O(n^2 \log n)$ .

### 3. A Greedy Algorithm for SCS

The greedy algorithm for the longest path problem can be restated for SCS as follows. Given a non-empty set of strings  $S$ , repeat the following step until  $S$  contains just one string.

Select a pair of strings  $s_1, s_2 \in S$  that maximizes  $\psi(s_1, s_2)$ . Remove  $s_1$  and  $s_2$  from  $S$ , replacing them with  $s_1 \circ s_2$ .

We refer to this algorithm as SGREEDY. Gallant [5] claims that  $\phi_{\text{SGREEDY}}(S) \leq (3/2)\phi^*(S)$ . This is not in fact true, as can be seen by considering the set of strings

$$S = \{\text{abc}bc\text{bc}bc\text{bc}, \text{c}bc\text{bc}bc\text{bc}, \text{bc}bc\text{bc}bc\text{bcd}\}$$

for which  $\phi_{\text{SGREEDY}}(S) = 20 > (3/2)\phi^*(S) = 19.5$ . One can easily generalize this example to show that there is no constant  $c < 2$  for which  $\phi_{\text{SGREEDY}}(S) \leq c\phi^*(S)$ . We currently do not know if there is some constant  $c > 2$  for which  $\phi_{\text{SGREEDY}}(S) \leq c\phi^*(S)$ .

On the other hand, Theorem 4 allows us to conclude that  $\psi^*(S) \leq 3\psi_{\text{SGREEDY}}(S)$ . In fact, we can improve the constant factor to 2 by noting that the instances of LPP that arise from the transformation from SCS have a special structure which is described in the following lemma.

**LEMMA 1.** *Let  $S$  be any set of strings and let  $(G, \ell) = \text{LPP}(S)$ . If  $\{w, x, y, z\} \subseteq V$  with  $\ell(w, y) = \max\{\ell(w, y), \ell(w, z), \ell(x, y), \ell(x, z)\}$  then  $\ell(w, y) + \ell(x, z) \geq \ell(w, z) + \ell(x, y)$ .*

*proof.* Identify  $w, x, y, z$  with the corresponding strings in  $S$  and let  $\alpha = \ell(w, y)$ ,  $\beta = \ell(x, z)$ ,  $\gamma = \ell(w, z)$ ,  $\delta = \ell(x, y)$ . Note that if  $\alpha \geq \gamma + \delta$  the result follows immediately. We will assume therefore that  $\alpha \leq \gamma + \delta$ .

We define some notation for designating substrings. If  $s = a_1 \dots a_r$  is a string,  $s[i]$  denotes the symbol  $a_i$  if  $i > 0$  and  $a_{r+i+1}$  if  $i < 0$ . The notation  $s[i, j]$  denotes the substring  $s[i] \dots s[j]$ .

By definition of  $G$ ,

$$w[-\alpha, -1] = y[1, \alpha] \quad (1)$$

$$w[-\gamma, -1] = z[1, \gamma] \quad (2)$$

$$x[-\delta, -1] = y[1, \delta] \quad (3)$$

From this we find

$$\begin{aligned} z[1, \gamma + \delta - \alpha] &= w[-\gamma, \delta - \alpha - 1] && \text{from (2) and } \alpha \leq \gamma + \delta \\ &= y[\alpha - \gamma + 1, \delta] && \text{from (1)} \\ &= x[\alpha - \gamma - \delta, -1] && \text{from (3)} \end{aligned}$$

Hence,  $\beta = \psi(x, z) \geq \gamma + \delta - \alpha$ .  $\square$

**THEOREM 5.** *Let  $S$  be any set of strings.  $\psi^*(S) \leq 2\psi_{\text{SGREEDY}}(S)$ .*

*Proof.* Let  $(G, \ell) = \text{LPP}(S)$ . Let  $H = \{h_1, \dots, h_s\}$  be the set of edges chosen by the greedy algorithm in the order in which they were selected (that is,  $h_1$  was selected first,  $h_2$  second, and so forth). Define  $\lambda_i^*$  to be the length of a longest path set on the subgraph of  $G$  defined by the set of edges that are still permissible after the selection of  $h_i$ . We show that for  $i \in [1, s]$ ,  $\lambda_{i-1}^* \leq 2\ell(h_i) + \lambda_i^*$ . This in turn, implies the theorem.

Let  $h_i = (w, y)$  and let  $X$  be a longest path in the subgraph of  $G$  defined by the set of edges that are permissible just before the selection of  $h_i$  (so  $\ell(X) = \lambda_{i-1}^*$ ). By definition of the greedy algorithm,  $\ell(w, y) = \max\{\ell(e) \mid e \in X\}$ . At most three edges of  $X$  are not permissible after  $h_i$  is selected. If at most two become impermissible, then clearly  $\lambda_i^* \geq \lambda_{i-1}^* - 2\ell(w, y)$  as desired.

If three edges become impermissible then one must have the form  $(w, z)$  with  $z \neq y$ , another the form  $(x, y)$  with  $x \neq w$  and the third one,  $e$  joins the last vertex on the path containing  $y$  with the first vertex on the path containing  $w$ . This means that  $X \cup \{h_1, \dots, h_{i-1}\}$  contains a path from  $x$  to  $z$ , which means that  $(x, z)$  is permissible *after*  $(w, y)$  is selected. By Lemma 1,  $\ell(w, y) + \ell(x, z) \geq \ell(w, z) + \ell(x, y)$ . Consequently,

$$2\ell(w, y) + \ell(X \cup \{(x, z)\} - \{e, (x, y), (w, z)\}) \geq \ell(X)$$

which implies  $\lambda_i^* \geq \lambda_{i-1}^* - 2\ell(w, y)$ .  $\square$

A naive implementation of SGREEDY takes at least quadratic time. A similar running time is obtained if one uses the transformation to LPP and then uses PGREEDY. A much faster running time can be obtained however by making use of an appropriate data structure. For our purposes, a *suffix tree*  $T$  is an abstract data type representing a collection of strings  $S = \{s_1, \dots, s_n\}$  on which the following operations are defined.

$T.\text{SUFFIX\_TREE}(S = \{s_1, \dots, s_n\})$  Initialize  $T$  to represent the strings in  $S$ . This operation may only be performed once.

$T.\text{LOOKUP}(\text{integer } i, j)$  Returns a pair  $[\ell, k]$ , where  $\ell$  is the length of the longest prefix of  $s_i$  which is also a suffix of some string in  $S - \{s_j\}$  and  $s_k$  is one such string.

$T.\text{DELETE}(\text{integer } i)$  Removes  $s_i$  from the set of strings represented.

The obvious implementation of a suffix tree is a *trie* (see [1]) containing an entry for every suffix of every string in the set. We assume that in addition to the edge labels, the nodes contain lists of integers, identifying the strings in the set that have suffixes ending at that node. This representation doesn't quite satisfy our needs, as the size of the trie and hence the time required to construct it is  $\Omega(m^2)$  in the worst-case. A more compact representation can be obtained by labeling edges with strings rather than single characters. This yields a representation that requires  $O(m)$  space and that can be constructed in  $O(m)$  time, as described by McCreight [11]. See also [2,12]. (Actually, McCreight defines a suffix tree to contain suffixes of a single string rather than a collection of strings. Our variant requires only minor modifications to McCreight's method.)

We perform deletion in suffix trees using lazy deletion. That is, to delete a string  $s_i$ , we simply mark it deleted in an auxiliary bit vector maintained for this purpose. When a lookup operation is performed, we probe the tree to find the longest match. Let  $u$  be the node at which the probe terminates. The list of matching strings at  $u$  is scanned and any that are marked deleted are removed from the list. If this makes the list empty and  $u$  has no children, then  $u$  is removed from the tree. If no acceptable match can be found in the list, the search continues at the parent of  $u$ .

The time required for a single lookup operation may well exceed the length of the string being searched for. However, any excess time is spent deleting node list entries. Since there are initially  $m+n$  such entries in the whole tree, the time spent on any sequence of lookups is  $O(m)$  plus the sum of the lengths of all strings being searched for.

We say that a sequence of lookup and delete operations is *monotonic* if for every  $i, j, k$  with  $j \neq k$ , whenever the sequence contains the operation  $T.LOOKUP(i, j)$  and later on the operation  $T.LOOKUP(i, k)$ , it contains the operation  $T.DELETE(j)$  between the other two.

We can speed up a monotonic sequence of operations, by maintaining for each string, a pointer to the node where the most recent lookup for that string ended. This allows us to avoid the initial probe of the tree when we perform a lookup operation. Instead, we use the pointer to go straight to the node where the last probe ended, and search up from that node if necessary. In this way, we can perform a monotonic sequence of  $r$  operations in  $O(m+r)$  time.

This analysis assumes that the symbol alphabet is small enough that it is reasonable to use a vector of pointers to children in each node, indexed by the first symbol of the strings labeling the edges. If a large alphabet is needed, a hash table may be used. Another option is to use a variant on Sleator and Tarjan's *lexicographic splay tree* [14]. With this representation, the time required to perform a sequence of operations is  $O(m)$ , plus the sum of the lengths of the strings being searched for, plus  $O(\log m)$  per operation. For a monotonic sequence of  $r$  operations the time is  $O(m+r \log m)$ .

An efficient implementation of the greedy algorithm for strings is shown in Figure 5. The algorithm does not explicitly combine strings, but keeps track of the decisions made using the two mappings  $left(i)$  and  $right(i)$ , which give the left and right neighbors of string  $i$  in the solution constructed so far. A value of 0 means that there is no neighbor. The solution is returned in these mappings. If a string  $i$  has no left neighbor yet,  $rightend(i)$  is the original string which is currently rightmost in the piece of the partial solution that contains  $i$ ;  $leftend(j)$  is similar.

The heap  $h$ , is used to determine which pair of strings should be combined next. Each string is entered in  $h$  with the *key* being the length of the best match for  $h$ . As the algorithm proceeds, certain matches become unavailable and the values of  $key()$  may become invalid. Consequently, whenever a string  $s_i$  is selected from  $h$ , a new lookup operation is performed in  $T$ . If the result of that operation is a match of the same length as  $key(i)$ , the strings are combined. If the lookup results in a shorter match, the value of  $key(i)$  is changed and the position of  $s_i$  in the heap is adjusted to reflect the new value. Note that a string is deleted from the heap once it is successfully matched with another string on its left end. Similarly, a string is deleted from the suffix tree once it is matched with a string on the right.

The running time of the algorithm is dominated by the operations on the data structures within the main loop. The number of iterations of the main loop is  $O(m)$  in the worst case. Since the heap operations are  $O(\log n)$  per operation, using the usual implementation, the total time spent on the heap operations is  $O(m \log n)$ . This can be improved to  $O(m+n \log n)$  by using Fibonacci heaps [3]. Since the sequence of operations on the suffix tree is monotonic, the time needed for the suffix tree operations is  $O(m)$  assuming a small alphabet and  $O(m \log m)$  assuming a large alphabet and the use of lexicographic splay trees. This yields a total running time of  $O(m+n \log n)$  for small alphabets and  $O(m \log m)$  for large alphabets.

```

procedure SGREEDY(string_set  $S = \{s_1, \dots, s_n\}$ , mapping  $left, right : [1, n] \mapsto [0, n]$ )
  integer  $i, j, \ell$ ; mapping  $key : [1, n] \mapsto \text{integer}$ ;
  mapping  $leftend, rightend : [1, n] \mapsto [1, n]$ ;
  heap  $h$ ; suffix_tree  $T$ ;
   $T.SUFFIX\_TREE(S)$ ;
  for  $i \in [1, n] \rightarrow$ 
     $left(i), right(i) := 0$ ;
     $leftend(i), rightend(i) := i$ ;
     $[key(i), j] := T.LOOKUP(i, i)$ ;
     $h.INSERT(i)$ ;
  rof;
  do  $|h| > 1 \rightarrow$ 
     $i := h.FINDMAX()$ ;
     $[\ell, j] := T.LOOKUP(i, rightend(i))$ ;
    if  $\ell = key(i) \rightarrow$ 
       $left(i), right(j) := j, i$ ;
       $leftend(rightend(i)) := leftend(j)$ ;
       $rightend(leftend(j)) := rightend(i)$ ;
       $T.DELETE(j); h.DELETE(i)$ ;
    |  $\ell < key(i) \rightarrow$ 
       $key(i) := \ell; h.SIFTDOWN(i)$ ;
    fi;
  od;
end

```

Figure 5: Greedy Algorithm for SCS

#### 4. SCS and the Traveling Salesman Problem

In this section we relate SCS to the path version of the *traveling salesman problem* (TSP). An instance of the traveling salesman problem is a list of cities  $C = (c_1, \dots, c_n)$  with a distance  $d(c_i, c_j)$  between each pair of cities. The object is to find a permutation  $\pi$  on  $\{1, \dots, n\}$ . that minimizes

$$\theta_\pi(C, d) = \sum_{i=1}^{n-1} d(c_{\pi_i}, c_{\pi_{i+1}})$$

We define  $\theta^*(C, d) = \min_\pi \theta_\pi(C, d)$ .

Let  $S = \{s_1, \dots, s_n\}$  be an instance of SCS. We define  $TSP(s_1, \dots, s_n)$  to be an instance  $(C, d)$  of TSP with  $C = (c_1, \dots, c_n, c_{n+1})$  and

$$d(c_i, c_j) = \begin{cases} |s_i| - \psi(s_i, s_j) & 1 \leq i \leq n \quad 1 \leq j \leq n \quad i \neq j \\ |s_i| & 1 \leq i \leq n \quad j = n+1 \\ \|S\| & i = n+1 \quad 1 \leq j \leq n \end{cases}$$

Note that in general, if  $\pi$  satisfies  $\theta_\pi(C, d) = \theta^*(C, d)$  then  $\pi_{n+1} = n+1$ .

**THEOREM 6.** *Let  $S = \{s_1, \dots, s_n\}$  be an instance of SCS,  $(C, d) = TSP(S)$  and let  $\pi$  be a permutation on  $\{1, \dots, n, n+1\}$  for which  $\pi_{n+1} = n+1$ . Then  $\theta_\pi(C, d) = \phi_{\pi'}(S)$ , where  $\pi'$  is the restriction of  $\pi$  to  $\{1, \dots, n\}$ . In particular,  $\theta^*(C, d) = \phi^*(S)$ .*

Theorem 6 implies that any good approximation algorithm for this version of the traveling salesman problem is a good approximation algorithm for SCS as well. The particular instances of TSP constructed by the transformation defined above have some special properties. First, they may be asymmetric; that is  $d(c_i, c_j)$  need not equal  $d(c_j, c_i)$ . The next theorem shows that they obey the so-called *triangle inequality*.

**THEOREM 7.** *Let  $S = \{s_1, \dots, s_n\}$  be an instance of SCS and let  $(C, d) = TSP(S)$ . For all  $c_i, c_j, c_k \in C$ ,  $d(c_i, c_k) \leq d(c_i, c_j) + d(c_j, c_k)$ .*

There exist efficient approximation algorithms for the symmetric version of the traveling salesman problem with triangle inequality that produce solutions within a factor of  $(3/2)$  of optimal. When the triangle inequality doesn't hold, finding good approximate solution is as difficult as finding optimum solutions (see [8]). For the asymmetric version with triangle inequality however, little is known. There are no known approximation algorithms that are both efficient and have good worst-case performance (nor have we found any), but the approximation problem has not been shown to be hard. Consequently, the relationship between SCS and TSP has yet to yield any immediately useful results. The relationship does imply some consequences if the status of either problem is resolved in the future. If good approximation algorithms are found for TSP, they may be applied to SCS. If the approximation problem for SCS is shown to be hard, then the approximation problem for TSP must be hard. It also may be that good approximation algorithms discovered for SCS, could be adapted to TSP, although this does not necessarily follow. Finally, a proof that TSP is hard to approximate would not imply that SCS is hard to approximate, but it would imply that approximation algorithms for SCS would have to exploit special structural properties present in the instances of TSP that result from our transformation. We close by mentioning one such property that may prove useful.

LEMMA 2. *Let  $S$  be any set of strings and let  $(C, d) = \text{TSP}(S)$ . If  $\{w, x, y, z\} \subset C$  with  $d(w, y) = \min\{d(w, y), d(w, z), d(x, y), d(x, z)\}$  then  $d(w, y) + d(x, z) \leq d(w, z) + d(x, y)$ .*

We refer to this as the *quadrilateral inequality* and note that it is analogous to the property described in Lemma 1.

## 5. Closing Remarks

We have related the shortest common matching string problem to two other problems, using transformations that preserve solution values. These two transformations reflect different ways of viewing SCS. We have used the transformations to gain insight into the problem of approximating SCS and have discovered several algorithms that have provably good performance with respect to the overlap measure. The best of these is the string version of the greedy algorithm for which we have described an efficient implementation using suffix trees.

While we have shown that the string version of the greedy algorithm has good worst-case performance with respect to the overlap measure, we cannot determine its performance with respect to the length measure. We know that it can be off by as much as a factor of two with respect to the length measure, but we don't know if it can be worse than this. One open problem then, is to resolve this issue.

Although, we have been unable to make use of the relationship between SCS and TSP to advantage, we feel that it may yet prove useful. More generally, we think that the use of transformations that preserve solution values can be used to extend the application of known approximation algorithms to new domains. A methodical development of such transformations could provide many useful results.

Another worthwhile line of investigation is to study the probable performance of the various approximation algorithms using appropriate probability models. It appears likely for example, that the directed matching algorithm for the longest path problem performs much better than its worst-case bound would indicate for a wide class of natural probability models. Similarly, one would expect the greedy algorithm to perform well in a probabilistic sense for many useful probability models.

## References

- [1] Aho, Alfred V., John E. Hopcroft and Jeffery D. Ullman. "Data Structures and Algorithms," Addison-Wesley, 1982.
- [2] Chen, M. T. and J. I. Seiferas. "Efficient and Elegant Subword-Tree Construction," University of Rochester, Computer Science Department, technical report TR 129, 12/83.
- [3] Fredman, Michael L. and Robert E. Tarjan. "Fibonacci Heaps and Their Uses in Improved Network Optimization Algorithms," *Proceedings of the IEEE Symposium on the Foundations of Computer Science*, 10/84, 338-346.
- [4] Gallant, John K., David Maier, James A. Storer. "On Finding Minimal Length Superstrings," *Journal of Computer and System Sciences*, vol. 20, no. 1, 50-58, 2/80.

- [5] Gallant, John K. "String Compression Algorithms," Ph.D. Dissertation, Princeton University, Department of Electrical Engineering and Computer Science, June 1982.
- [6] Gingeras, T. R., J. P. Milao, D. Sciaky and R. J. Roberts. "Computer Programs for the Assembly of DNA Sequences," *Nucleic Acids Research*, vol. 7, 1979, 529-545.
- [7] Karp, Richard M. "A Patching Algorithm for the Nonsymmetric Traveling Salesman Problem," *SIAM Journal on Computing*, vol. 8, no. 4, 561-573, 11/79.
- [8] Lawler, E. L., J. K. Lenstra and A. H. G. Rinnooy Kan (editors). "The Traveling Salesman Problem," John Wiley and Sons, Ltd, 1986.
- [9] Maier, David and James A. Storer. "A Note on the Complexity of the Superstring Problem," Princeton University Technical Report 233, Department of Electrical Engineering and Computer Science, October 1977.
- [10] Mayne, A. and E. B. James. "Information Compression by Factorising Common Strings," *The Computer Journal* 18, 1975, 157-160.
- [11] McCreight, Edward M. "A Space-Economical Suffix Tree Construction Algorithm," *Journal of the ACM*, vol. 23, 4/76, 262-272.
- [12] Rodeh, M., V. R. Pratt and S. Even. "Linear Algorithm for Data Compression via String Matching," *Journal of the ACM*, vol. 28, 1/81, 16-24.
- [13] Shapiro, M. B. "An Algorithm for Reconstructing Protein and RNA Sequences," *Journal of the ACM*, vol. 14, 1967, 720-731.
- [14] Sleator, Daniel D. and Robert E. Tarjan. "Self-Adjusting Binary Search Trees," *Journal of the ACM*, vol. 32, 7/85, 652-686.
- [15] Stefik, M. "Inferring DNA Structures From Segmentation Data," *Artificial Intelligence*, vol. 11, 1978, 85-114.
- [16] Storer, James A. and Thomas G. Szymanski. "Data Compression via Textual Substitution," *Journal of the ACM*, vol. 29, 10/82, 928-951.
- [17] Tarjan, Robert E. "Data Structures and Network Algorithms," Society for Industrial and Applied Mathematics, 1983.